

APPLICATION
FOR
UNITED STATES LETTERS PATENT

TITLE: MICROINSTRUCTION POINTER STACK IN A
PROCESSOR

APPLICANT: MICHAEL P. CORNABY AND BEN CHAFFIN

Scott C. Harris
Fish & Richardson P.C.
4350 La Jolla Village Drive
Suite 500
San Diego, CA 92122
Telephone: 858-678-5070
Facsimile: 858-678-5099

CERTIFICATE OF MAILING BY EXPRESS MAIL

Express Mail Label No. EU 047 039 083 US

I hereby certify that this correspondence is being deposited with the United States Postal Service as Express Mail Post Office to Addressee with sufficient postage on the date indicated below and is addressed to the Commissioner for Patents, Washington, D.C. 20231.

12-20-01
Date of Deposit

Gabe Lewis
Signature

Gabe Lewis
Typed or Printed Name of Person Signing
Certificate

MICROINSTRUCTION POINTER STACK IN A PROCESSOR

TECHNICAL FIELD

This invention relates to a microinstruction pointer stack in a processor.

BACKGROUND

5 A microprocessor is a computer processor on a microchip. The microprocessor is typically designed to perform arithmetic and logic operations that make use of small number-holding areas called registers. Typical microprocessor operations include adding, subtracting, comparing, and fetching operands from
10 memory or registers. These operations result from execution a set of instructions that comprise a program. The set of instructions are part of the microprocessor design.

DESCRIPTION OF DRAWINGS

FIG. 1 is a block diagram of a processor.

15 FIG. 2 is a block diagram of an executive environment of the processor of FIG. 1.

FIG. 3 is a diagram of an out of order microinstruction pointer stack implemented in the out of order execution core of FIG. 1.

20 FIG. 4 is a flow diagram of a μ IP stack process.

DETAILED DESCRIPTION

Referring to FIG. 1 a processor 10 is shown. The processor 10 is a three way super scaler, pipelined architecture. Using parallel processing techniques, the processor 10 decodes, dispatches, and completes execution of (retire) three instructions per clock cycle. To handle this level of instruction throughput, the processor 10 uses a decoupled, e.g., twelve stage pipeline that supports out of order instruction execution. The pipeline of the processor 10 is divided into four sections, i.e., a first level cache 12, a second level cache 14, a front end 16, an out of order execution core 18, and a retire section 20. Instructions and data are supplied to these units through a bus interface unit 22 that interfaces with a system bus 24. The front end 16 supplies instructions in program order to the out of order execution core 18 that has very high execution bandwidth and can execute basic integer operations with one-half clock cycle latency. The front end 16 fetches and decodes instructions into simple operations called micro-ops (μ ops). The front end 16 can issue multiple μ ops per cycle, in original program order, to the out of order execution core 18. The front end 16 performs several basic functions. For example, the front end 16 performs prefetch instructions that are likely to be executed. The front end 16 decodes instructions into micro operations and generates micro code for

complex instructions, delivers decoded instructions from an execution trace cache 26, and predicts branches using advanced algorithms in a branch prediction unit 28.

The front end 16 of the processor 10 addresses some common problems in high speed, pipelined microprocessors. Two of these problems, for example, contribute to major sources of delays. These problems are the time to decode instructions fetched from the target and time wasted to decode instructions due to branches or branch targets that occur in the middle of cache lines.

The execution trace cache 26 addresses both of these issues by storing decoded instructions. Instructions are fetched and decoded by a translation engine (not shown) and built into sequences of μops called traces. These traces of μops are stored in the trace cache 26. The instructions from the most likely target of a branch immediately follow the branch, without regard for continuity of instruction addresses. Once a trace is built, the trace cache 26 is searched for the instruction that follows that trace. If that instruction appears as the first instruction in an existing trace, the fetch and decode of instructions 30 from the memory hierarchy ceases and the trace cache 26 becomes the new source of instructions.

The execution trace cache 18 and the translation engine (not shown) have cooperating branch prediction hardware. Branch

targets are predicted based on their linear addresses using Branch Target Buffers (BTBS) 28 and fetched as soon as possible. The branch targets are fetched from the trace cache 26 if they are indeed cached there; otherwise, they are fetched from the memory hierarchy. The translation engine's branch prediction information is used to form traces along the most likely paths.

The core 18 executes instructions out of order enabling the processor 10 to reorder instructions so that if one μ op is delayed while waiting for data or a contended execution resource, other μ ops that are later in a program order may execute before the delayed μ ops. The processor 10 employs several buffers to smooth the flow of μ ops. This implies that when one portion of the pipeline experiences a delay, that delay may be covered by other operations executing in parallel or by the execution of μ ops which were previously queued up in one of the buffers.

The core 18 is designed to facilitate parallel execution. The core 18 can dispatch up to six μ ops per cycle; note that this exceeds the trace cache 26 and retirement 20 μ op bandwidth. Most pipelines can start executing a new μ op every cycle, so that several instructions can be processed any time for each pipeline. A number of arithmetic logical unit (ALU) instructions can start two per cycle, and many floating point instructions can start one every two cycles. Finally, μ ops can

begin execution, out of order, as soon as their data inputs are ready and resources are available.

The out of order execution core 18 includes an out of order microinstruction pointer (IP) stack 100. In general, a stack is
5 a data area or buffer used for storing requests that need to be handled. A stack is typically a push-down list, meaning that as new requests come into the stack, the stack pushes down older requests. Another way of looking at a push-down list - or stack - is that a program usually takes its next item to handle from
10 the top of the stack, unlike other arrangements such as "FIFO" or "first-in first-out" buffers. The stack 100 is implemented in a microcode environment. This allows fast subroutine returns in microcode. It also allows fast assist returns in microcode.

The μ IP stack 100 is different from a macroinstruction
15 stack in several ways. For example, the μ IP stack 100 is not visible from a system level (i.e., the μ IP stack 100 cannot be directly manipulated from macrocode). The μ IP stack 100 is an out-of-order stack where values are placed on the stack and removed from the stack before it is known if the sequence of
20 operations were valid. Thus, a set of recovery mechanisms is required to correct a sequence of operations that are later determined to be incorrect. The process of manipulating the stack (PUSH, POP, etc.) is very different from a traditional

macroinstruction stack because of the out-of-order nature of the stack 100.

The μ IP stack 100 provides a mechanism for improving the performance of microcode (μ code) execution. Microcode is programming that is ordinarily not program-addressable but, unlike hardwired logic, is capable of being modified. Microcode may sometimes be installed or modified by a device's user by altering programmable read-only memory (PROM) or erasable programmable read-only memory (EPROM).

The μ IP stack 100 provides a lower-overhead ability to jump to various subroutines and use "assists" to efficiently accomplish μ code functions. The μ IP stack 100 has significant performance and μ code efficiency implications that permeate numerous instructions. For example, use of the μ IP stack 100 improves performance by removing indirect μ code jumps and allows μ code to share routines more easily by removing subroutine penalties. By removing subroutine penalties, the μ IP stack 100 allows μ code to share routines more easily.

The retirement section 20 receives the results of the executed μ ops from the execution core 18 and processes the results so that the proper architectural state is updated according to the original program order. For semantically correct execution, the results of instructions are committed in original program order before the instructions are retired.

Exceptions may be raised as instructions are retired. Thus, exceptions do not occur speculatively, but rather exceptions occur in the correct order, and the processor 10 can be correctly restarted after execution.

5 When a μ op completes and writes its result to the destination, it is retired. Up to three μ ops may be retired per cycle. A ReOrder Buffer (ROB) (not shown) in the retirement section 20 is the unit in the processor 10 which buffers completed μ ops, updates the architectural state in order, and
10 manages the ordering of exceptions.

 The retirement section 20 also keeps track of branches and sends updated branch target information to the BTB 28 to update branch history. In this manner, traces that are no longer needed can be purged from the trace cache 26 and new branch
15 paths can be fetched, based on updated branch history information.

 Referring to FIG. 2, an execution environment 50 is shown. Any program or task running on the processor 10 (of FIG. 1) is given a set of resources for executing instructions and for
20 storing code, data, and state information. These resources make up the execution environment 50 for the processor 10. Application programs and the operating system or executive running on the processor 10 use the execution environment 50 jointly. The execution environment 50 includes basic program

execution registers 52, an address space 54, Floating Point Unit (FPU) registers 56, multimedia extension registers (MMX) 58, and SIMD extension registers 60.

Any task or program running on the processor 10 can address a linear address base 54 of up to four gigabytes (2^{32} bytes) and a physical address space of up to 64 gigabytes (2^{36} bytes). The address space 54 can be flat or segmented. Using a physical address extension mechanism, a physical address space of 2^{36-1} can be addressed.

The basic program execution registers 52 include eight general purpose registers 62, six segment registers 64, an EFLAGS register 66, and an EIP (instruction pointer) register 68. The basic program execution registers 52 provide a basic execution environment in which to execute a set of general purpose instructions. These instructions perform basic integer arithmetic on byte, word, and doubleword integers, handle program flow control, operate on bit and byte strengths, and address memory.

The FPU registers 56 include eight FPU data registers 70, an FPU control register 72, a status register 74, an FPU instruction pointer register 76, an FPU operand (data) pointer register 78, an FPU tag register 80 and an FPU op code register 82. The FPU registers 56 provide an execution environment for operating on single precision, double precision, and double

extended precision floating point values, word-, doubleword, and quadword integers, and binary coded decimal (BCD) values.

The eight multimedia extension registers 58 support execution of single instruction, multiple data (SIMD) operations
 5 on 64-bit packed byte, word, and doubleword integers.

The SIMD extension registers 60 include eight extended multimedia (XMM) data registers 84 and an MXCSR register 86. The SIMD extension registers 60 support execution of SIMD operations on 128-bit packed single precision and double
 10 precision floating point values and on 128-bit packed byte, word, doubleword and quadword integers.

A stack (not shown) supports procedure or subroutine calls and the passing of parameters between procedures or subroutines.

The general purpose registers 62 are available for storing
 15 operands and pointers. The segment registers 64 hold up to six segment selectors. The EFLAGS (program status and control) registers 66 report on the status of a program being executed and allows limited (application program level) control of the processor. The EIP (instruction pointer) register 68 has a 32-
 20 bit pointer to the next instruction to be executed.

The 32-bit general purpose registers 62 are provided for holding operands for logical and arithmetic operations, operands for address calculations, and memory pointers. The segment registers 64 hold 16-bit segment selectors. A segment selector

is a special pointer that identifies a segment in memory. To access a particular segment in memory, the segment selector for that segment must be present in the appropriate segment register 64.

5 When writing application code, programmers generally produce segment selectors with assembler directives and symbols. The assembler and other tools generate the actual segment selector values associated with these directives and symbols. If writing system code, programmers may need to generate segment
10 selectors directly.

How segment registers 64 are used depends on the type of memory management model that the operating system or executive is using. When using a flat (unsegmented) memory model, the segment registers 64 are loaded with segment selectors that
15 point to overlapping segments, each of which begins at address zero on the linear address space. These overlapping segments also include the linear address space for the program. Typically, two overlapping segments are defined: one for code and another for data and stacks. The CS segment register (not
20 shown) of the segment registers 64 points to the code segment and all other segment registers point to the data and stack segment.

When using a segmented memory model, each segment register 64 is ordinarily loaded with a different segment selector so

that each segment register 64 points to a different segment within the linear address space. At any time, a program can thus access up to six segments in the linear address space. To access a segment not pointed to by one of the segment registers 5 64, a program first loads the segment selector to be accessed into a segment register 64.

The 32-bit EFLAGS register 66 has a group of status flags, a control flag, and a group of system flags. Some of the flags in the EFLAGS register 66 can be modified directly, using 10 special purpose instructions. The following instructions can be used to move groups of flags to and from the procedure stacks or general purpose register: LAHF, SAHF, push-F, push-FD, pop-F, and pop-FD. After the contents of EFLAGS register 66 have been transferred to the procedure stack or a general purpose 15 register, the flags can be examined and modified using the processor 10 bit manipulation instructions.

When suspending a task, the processor 10 automatically saves the state of the EFLAGS register 66 in the task state segment (TSS) (not shown) for the task being suspended. When 20 binding itself to a new task, the processor 10 loads the EFLAGS register 66 with data from the new tasks program state register (PSS, not shown).

When a call is made to an interrupt or an exception handler procedure the processor 10 automatically saves the state of the

EFLAGS register 66 on the procedure stack. When an interrupt or exception is handled with a task switch, the state of the EFLAGS register 66 is saved on the TSS for the task being suspended.

The fundamental data types used in the processor 10 are
 5 bytes, words, doublewords, quadwords and double quadwords. A byte is eight bits, a word is two bytes (16-bits), a doubleword is four bytes (32-bits), a quad word is eight bytes (64-bits), and a double quadword is sixteen bytes (128-bits).

Referring to FIG. 3, the first n entries of the μ IP stock
 10 100 are the *in flight* part of the μ IP stack 100. *In flight* entries refer to entries currently being processed. The other entries are the retired part of the μ IP stack. Retired entries are those that are no longer being processed.

A μ IP field 104 has the μ IP pushed by a `ms_push` μ Op
 15 (described below) and used by a `fast_return` μ Op (described below) and has a width of 14 bits.

A `BackPtr` field 106 points to a next entry in the μ IP stack
 for μ TOS to point to after an `ms_return/ms_pop` μ op. It has room
 for the pointer value and a wrap bit so its width depends on
 20 stack size.

When an *in flight* entry retires, the `RetPtr` field 102 is
 updated to point to the location in the retired stack (not
 shown) to which the entry is copied. Thus, its width depends on
 the stack size.

A RO/RI field 108 records whether this in flight entry has retired. Two bits are needed to handle wrap cases and thus its width is 2 bits.

5 The μ IP stack 100 includes four pointers that point to different entries in the μ IP stack 100. The four pointers are a μ TOS pointer 110, an μ Alloc pointer 112, a NextRet pointer 114, and a μ RetTOS pointer 116. The μ TOS pointer 110, μ Alloc pointer 112, and NextRet pointer 114 require a wrap bit.

10 The μ TOS pointer 110 is the current top of stack 100 for μ Op issue and points to one of the entries in the table or to a NULL entry. The μ TOS pointer 110 is set to the current μ Alloc pointer 112 on the issue of a `ms_push` μ Op (described below). Note that it can point to any entry in the table (both in the in flight section and the retired section).

15 The μ Alloc pointer 112 points at the next entry to be allocated when an `ms_push` μ Op (described below) is issued. The last entry this pointer can point to is $n-1$. After this point it wraps, so the entries from 0 to $n-1$ are treated as a circular queue.

20 The NextRet pointer 114 points at the next entry to be deallocated from the μ IP stack 100 when a μ IP stack operation retires. Like the μ Alloc pointer 112, this pointer wraps at $n-1$.

The μ RetTOS pointer 116 points at the retired top of stack 100. This pointer can never point to entries 0 to $n-1$.

Additional μ Ops are used with the μ IP stack 100. The additional μ Ops are: `ms_call`, `ms_push`, `ms_pop`, `ms_return`, `ms_tos_read`, and `ms_ μ ip_stack_clear`. Alternatively, `call`, `return`, and `clear` could be attached to other μ ops.

5 The `ms_call` μ OP takes the next μ ip, pushes it on the μ IP stack 100, and uses the μ ip in the immediate field as the target μ ip of a jump.

The `ms_push` μ OP takes the value in the immediate field and pushes it on the μ IP stack 100.

10 The `ms_pop` μ OP pops a value off the μ IP stack 100 and replaces this μ Op's immediate field.

The `ms_return` μ OP pops a value off the μ IP stack 100 and jumps to that μ ip.

15 The `ms_tos_read` μ OP reads a value off the μ IP stack 100 and replaces this μ Op's immediate field, without changing the contents of the μ IP stack 100.

The `ms_ μ ip_stack_clear` μ OP sets the μ IP stack pointers to the reset values. Note that this μ Op is executed at a time when all preceding stack operations and retirements are complete.

20 Referring to FIG. 4, a micro-instruction pointer (μ IP) stack process 200 includes executing (202) microcode (μ code) stored in a out-of-order μ IP stack. The process 200 pushes (204) a next μ IP on to the μ IP stack and uses the next μ IP in an intermediate field as a target μ IP in a jump operation. The

process 200 takes (206) a value of an intermediate field of a microoperation (μOp) and pushes the value on to the μIP stack.

The process 200 pops (208) a value off the μIP stack and replaces a current μOp intermediate field with the value. The

5 process 200 pops (210) a value off of the μIP stack and jumps to that value.

The process 200 reads (212) a value off the μIP stack and replaces a μOp 's intermediate field with the value. The process 200 sets (214) the μIP stack pointers to reset.

10 The following terminology is used throughout the description below. `MAX_INFLIGHT` refers to the maximum number of calls allowed to be alive in the processor at once. `MAX_STACK` refers to the deepest function nesting level allowed. `RET_OFFSET` refers to offset in the table 100 of the first entry in the
15 retired area. `NULL_INDEX` refers to the index in the table 100 of the null entry. The code below assumes that this lies between the in flight section and the retired section of the stack 100.

At reset:

```

20       $\mu\text{TOS}.\text{ptr} = \text{NULL\_INDEX}$ 
       $\mu\text{TOS}.\text{wrap} = 0$ 
       $\mu\text{Alloc}.\text{ptr} = 0$ 
       $\mu\text{Alloc}.\text{wrap} = 0$ 
       $\text{NextRet}.\text{ptr} = 0$ 
       $\text{NextRet}.\text{wrap} = 0$ 
25       $\mu\text{RetTOS} = \text{NULL\_INDEX}$ 

```


On issue of ms_call μ Op:

```

    if ( $\mu$ Alloc.ptr == NextRet.ptr &&  $\mu$ Alloc.wrap !=
    NextRet.wrap) MSStall;
    stack[ $\mu$ Alloc.ptr].BackPtr =  $\mu$ TOS;
5    stack[ $\mu$ Alloc.ptr]. $\mu$ ip = current_ $\mu$ ip + 1;
    stack[ $\mu$ Alloc.ptr].R[ $\mu$ Alloc.wrap] = 0;
     $\mu$ TOS =  $\mu$ Alloc; //copies both the pointer and the wrap
    bit
     $\mu$ Alloc.ptr = ( $\mu$ Alloc.ptr + 1)%MAX_INFLIGHT;
10    if ( $\mu$ Alloc.ptr == 0)
         $\mu$ Alloc.wrap ^= 1;
    next_ $\mu$ ip = ms_call  $\mu$ ip (immediate field)

```

where, if the μ Alloc pointer is equal to the NextRet pointer and
15 their wrap bits are different, then the in flight table is full
and one cannot issue any more call/push μ ops until one retires.
If the table is not full, then μ Alloc.ptr points to the next
entry to be allocated, so it is updated. More specifically, the
current value of μ TOS is placed into the BackPtr so we know
20 where to return to. The μ IP of the μ op after the call/push is
put into the μ ip field. One of the R (retired) bits cleared
(which R bit one depends on the current wrap bit of μ Alloc). The
 μ TOS is set to point to the current entry (μ Alloc). Both the
pointer and the wrap bit must be copied. μ Alloc is incremented,
25 wrapping (and toggling the wrap bit) if necessary. Finally,
branch to the μ IP in the immediate field of the μ op.

On issue of ms_push μ Op instruction, the same events as in
a ms_call μ Op occur, except that the μ op's immediate field is
copied into the μ ip field of the stack instead of the μ IP of the

next μ op, and the next μ IP to be sequenced is determined as usual.

On issue of `ms_return μ Op` instruction:

```

5      next_μip = stack[μTOS.ptr].μip;
      back_ptr = stack[μTOS.ptr].BackPtr;
      if (stack[back_ptr.ptr].R[back_ptr.wrap] == 1)
      μTOS.ptr = stack[back_ptr.ptr].RetPtr; //wrap bit
      doesn't matter
10     else
      μTOS = back_ptr; //copies both pointer and wrap bit

```

where it gets the next μ IP to sequence from the μ ip field of the stack entry pointed to by μ TOS. Then pop the stack: the BackPtr of the entry pointed to by μ TOS has the index of the entry underneath this one on the stack. However, if that entry has retired since the BackPtr was set, it may have been overwritten by another speculatively issued call. So check the R bit of the entry pointed to by BackPtr. If it is 0, then the BackPtr entry is valid and we set μ TOS to point to it; if the R bit is 1, then the RetPtr field of that entry has its forwarding address.

On issue of the `ms_pop μ OP` instruction the same events occur as the `ms_return μ OP` instruction, except the immediate field of the `ms_pop μ op` is replaced with the μ ip field from the stack entry pointed to by μ TOS, and the next μ IP is determined normally.

On retirement of ms_call or ms_push μ Op instruction:

```

old_μRetTOS = μRetTOS;
μRetTOS++; //no wrap needed --better not overflow!
stack[μRetTOS].BackPtr.ptr = old_μRetTOS; //wrap bit
5  doesn't matter
  stack[μRetTOS].μip = stack[NextRet.ptr].μip;
  stack[NextRet.ptr].RetPtr = μRetTOS;
  stack[NextRet.ptr].R[NextRet.wrap]=1;
  if (NextRet.ptr == μTOS.ptr) //wrap bits always the
10  same
      μTOS.ptr = μRetTOS; //wrap bit doesn't matter
  NextRet.ptr = (NextRet.ptr +1)%MAX_INFLIGHT;
  if (NextRet.ptr == 0)
      NextRet.wrap ^=1;
15  clear any MSStall due to full in-flight stack;

```

where the μ RetTos is incremented, μ RetTOS, making sure it moves between NULL and the first entry correctly. The old value of the μ RetTOS is put in the BackPtr of the new retired entry. The μ IP from the entry pointed to by NextRet (the next entry to retire) is copied to the μ IP field of the new retired entry. The RetPtr of the entry pointed to by NextRet is set to the new μ RetTOS. The R bit of the entry pointed to by NextRet is set to 1. If the NextRet pointer equals the μ TOS pointer, then we have just

25 invalidated the entry pointed to by μ TOS, so set μ TOS to point to the retired copy (the new value of μ RetTOS). Increment NextRet, wrapping and toggling the wrap bit if necessary. Clear the MS stall condition resulting from too many push/call operations in flight.

30

On retirement of ms_return or ms_pop µOp instruction:

```
µRetTOS--;
-OR- µRetTOS = stack[µRetTOS].BackPtr;
```

5

where the µRetTOS pointer is decremented, or replaced with the BackPtr from the entry it points to; these are equivalent. The BackPtr is implemented for the retired stack since it is used in the manipulation of µTOS (unless the rule is: if µTOS is within the retired stack, decrement; otherwise follow the BackPtr).

10

On mispredicted macrobranch/microbranch:

```
µAlloc = mispred_µAlloc; //copies both pointer and
wrap bit if
(stack[mispred_µTOS.ptr].R[mispred_µTOS.wrap])
µTOS.ptr = stack[mispred_µTOS.ptr].RetPtr; //wrap
bit doesn't matter
else
µTOS = mispred_µTOS; //copies both pointer and wrap
bit
```

15

20

where the µAlloc and µTOS pointers are restored to the values that were saved when the branch which is mispredicting was issued. However, if the entry which the branch's µTOS points to has retired, set µTOS to point to its new location in the retired stack instead.

25

On trap or fault:

```
µTOS.ptr = NULL_INDEX
µTOS.wrap = 0
µAlloc.ptr = 0
```

30

```

    uAlloc.wrap = 0
    NextRet.ptr = 0
    NextRet.wrap = 0
    uRetTOS = NULL_INDEX

```

5

On assist:

```

    uRetTOS++;
    stack[uRetTOS].uip = assist uip
    stack[uRetTOS].BackPtr = uRetTOS - 1;
    uTOS = uRetTOS
    uAlloc.ptr = 0
    uAlloc.wrap = 0
    NextRet.ptr = 0
    NextRet.wrap = 0

```

10

15

In the case of a trap, the uIP stack 100 can be completely cleared. By definition of trap, all the previous flows are complete, and all the new flows are speculative, so all values on the uIP stack are speculative and can be thrown away.

20

There are two cases for a fault. If the fault will not return to the current flow of execution, the uIP stack 100 can be completely cleared. If the fault will return to the flow of execution, either the uIP stack 100 needs to be recovered or it needs to be cleared and a restriction placed on flows which can do this as to their use of ms_push/fast_return.

25

The following example illustrates operation of the uIp stack 100. Consider, for example, the following sequence of events occurring in the uIP stack 100:

(A) issue ms_call #1 from uip X

(B) issue ms_call #2 from uip Y

30

- (C) issue `μ_jump_cc #1` which will mispredict
- (D) issue `ms_ret` from call #2
- (E) issue `μ_jump_cc #2` which will mispredict
- (F) retire call #1
- (G) `μ_jump_cc #2` executes and mispredicts
- (H) `μ_jump_cc #1` executes and mispredicts
- (I) retire call #2

Below is the `μIP` stack 100 as it will appear after each of these operations, assuming `MAX_INFLIGHT=3` and `MAX_STACK=4`. The pointers are indicated on the right; the number to the right of the pointer is the wrap bit.

Start:

Entry	<code>μIP</code>	BackPtr	R0	R1	RetPtr
0					<code><μAlloc-0 <NextRet-0</code>
1					
2					
NULL	0	NULL	0	0	<code><μRetTOS <μTOS</code>
4			0	0	
5			0	0	
6			0	0	
7			0	0	

After (A): Push `X+1` onto the `μIP` stack 100, update `μAlloc` and `μTOS` pointers.

Entry	μ IP	BackPtr	R0	R1	RetPtr	
0	X+1	NULL	0			<NextRet-0 < μ TOS-0
1						< μ Alloc-0
2						
NULL	0	NULL	0	0		< μ RetTOS
4			0	0		
5			0	0		
6			0	0		
7			0	0		

After (B): Push Y +1 on the μ IP stack 100, update μ Alloc and μ TOS pointers.

Entry	μ IP	BackPtr	R0	R1	RetPtr	
0	X+1	NULL	0			<NextRet-0
1	Y+1	0-0	0			< μ TOS-0
2						< μ Alloc-0
NULL	0	NULL	0	0		< μ RetTOS
4			0	0		
5			0	0		
6			0	0		
7			0	0		

After (C): μ_jump_cc #1 issues, taking the values of μ Alloc=2-0 and μ TOS=1-0 with it.

After (D): Next μ IP is Y+1 (μ ip field of μ TOS entry). Take BackPtr of μ TOS entry (0-0): look up stack[BackPtr.ptr]

.R[Backptr.wrap]: stack[0].R0 indicates this entry has not retired and is still valid, so the μ TOS pointer gets 0-0.

Entry	μ IP	BackPtr	R0	R1	RetPtr	
5	0	X+1	NULL	0		<NextRet-0 < μ TOS-0
	1	Y+1	0-0	0		
	2					< μ Alloc-0
10	NULL	0	NULL	0	0	< μ RetTOS
	4			0	0	
15	5			0	0	
	6			0	0	
20	7			0	0	

After (E): $\mu_jump_cc_2$ issues, taking the values of μ Alloc=2-0 and μ TOS=0-0 with it.

After (F): Increment μ RetTOS and copy NextRet entry to new μ RetTOS entry. Set the RetPtr of the NextRet entry to point to its new location, and set the R bit. Increment NextRet.

Entry	μ IP	BackPtr	R0	R1	RetPtr	
30	0	X+1	NULL	1	4	< μ TOS-0
	1	Y+1	0-0	0		<NextRet-0
	2					< μ Alloc-0
35	NULL	0	NULL	0	0	
	4	X+1	NULL	0	0	< μ RetTOS
40	5			0	0	
	6			0	0	
45	7			0	0	

After (G): μ_jump_cc #2 mispredicts, returning $\mu Alloc=2-0$ and $\mu TOS=0-0$. Set $\mu Alloc$ to 2-0, no change. Check R bit of μTOS being restored--it is set, so set μTOS to the RetPtr of that entry instead.

5

Entry	μIP	BackPtr	R0	R1	RetPtr
0	X+1	NULL	1		4
1	Y+1	0-0	0		<NextRet-0
2					< $\mu Alloc$ -0
NULL	0	NULL	0	0	
4	X+1	NULL	0	0	< $\mu RetTOS$ < μTOS -0
5			0	0	
6			0	0	
7			0	0	

25

After (H): μ_jump_cc #1 mispredicts, returning $\mu Alloc=2-0$ and $\mu TOS=1-0$. Set $\mu Alloc$ to 2-0, no change. Check R bit of μTOS we're restoring -- it is not set, so set μTOS to the value returned by the mispredict.

30

Entry	μIP	BackPtr	R0	R1	RetPtr
0	X+1	NULL	1		4
1	Y+1	0-0	0		<NextRet-0 μTOS -0
2					< $\mu Alloc$ -0
NULL	0	NULL	0	0	
4	X+1	NULL	0	0	< $\mu RetTOS$ -0
5			0	0	

6				0	0	
7				0	0	

5

After (I): increment μRetTOS and copy NextRet entry to new μRetTOS entry. Set the RetPtr of the NextRet entry to point to its new location, and set the R bit. Since $\text{NextRet} == \mu\text{TOS}$, we have just retired the last valid entry on the μIP stack 100, so set μTOS to point to the new location of the current entry on the retired stack. Increment NextRet.

10

15

20

25

30

Entry	μIP	BackPtr	R0	R1	RetPtr	
0	X+1	NULL	1		4	
1	Y+1	0-0	0		5	<NextRet-0
2						< μAlloc -0 <NextRet-0
NULL	0	NULL	0	0		
4	X+1	NULL	0	0		
5	Y+1	4	0	0		< μRetTOS < μTOS -0
6			0	0		
7			0	0		

Several considerations can be made for debugging and design verification. For example, for patching considerations, the μRetTOS pointer can be readable and writeable through microcode. In addition, the retired instruction can be writeable through control register access. This allows microcode to clear the instruction from in flight stack. The microcode can thus read

the μ RetTOS to determine the number of entries on the retired stack and pop the entries off the stack 100. Popping entries off the stack 100 takes the entries to the EXECUTIVE where the entries can be examined. The microcode can restore the μ RetTOS
5 (which puts the stack back to the state it was before the pops), and modify the values in the μ RetTOS via control register writes.

The stack pointers μ TOS, μ Alloc, and NextRet should be visible for debugging. One way to make the stack pointers
10 viable is to allow access through a control register.

Access to the in flight μ IP stack 100 can be through a control register mechanism, but an array dump mechanism is acceptable.

Having control register access to the in flight μ IP stack
15 100 hardware may increase microcode flexibility at the risk of being extremely hard to maintain correctness.

A number of embodiments of the invention have been described. Nevertheless, it will be understood that various modifications may be made without departing from the spirit and
20 scope of the invention. For example, an option is to provide a path from the TBP μ to the MS where the EV_ μ IP can be accessed. This would allow the assisting μ Op's μ ip stack 100 to be pushed on the μ IP stack 100 and allow faster returns from assists. Alternately, another μ Op could be used to get the μ ip from the

EXEC to the MS for pushing on the MS stack. For longer assist flows, this could eliminate the indirect branch latency.

Accordingly, other embodiments are within the scope of the following claims.